Backoff and Interpolation

- Sometimes it helps to use **less** context
 - Condition on less context for contexts you haven't learned much about
- Backoff:
 - use trigram if you have good evidence,
 - otherwise bigram, otherwise unigram
- Interpolation:
 - mix unigram, bigram, trigram
- Interpolation works better

Linear Interpolation

• Simple interpolation

$$\hat{P}(w_n | w_{n-2} w_{n-1}) = \lambda_1(w_{n-2}^{n-1}) P(w_n | w_{n-2} w_{n-1}) \\
+ \lambda_2(w_{n-2}^{n-1}) P(w_n | w_{n-1}) \\
+ \lambda_3(w_{n-2}^{n-1}) P(w_n)$$

How to set the lambdas?

• Use a **held-out** corpus

- Choose λ s to maximize the probability of held-out data:
 - Fix the N-gram probabilities (on the training data)
 - Then search for λ s that give largest probability to held-out set:

$$\log P(w_1...w_n | M(/_1.../_k)) = \operatorname{alog} P_{M(/_1.../_k)}(w_i | w_{i-1})$$

Smoothing for Web-scale N-grams

- "Stupid backoff" (Brants et al. 2007)
- No discounting, just use relative frequencies

$$S(w_{i} | w_{i-k+1}^{i-1}) = \begin{cases} 1 \\ 1 \\ 1 \\ 1 \end{cases} \quad \frac{\text{count}(w_{i-k+1}^{i})}{\text{count}(w_{i-k+1}^{i-1})} & \text{if } \text{count}(w_{i-k+1}^{i}) > 0 \\ 1 \\ 1 \\ 1 \end{cases}$$

$$S(w_i) = \frac{\operatorname{count}(w_i)}{N}$$

Unknown words: Open versus closed vocabulary tasks

- If we know all the words in advanced
 - Vocabulary V is fixed
 - Closed vocabulary task
- Often we don't know this
 - Out Of Vocabulary = OOV words
 - Open vocabulary task
- Instead: create an unknown word token <UNK>
 - Training of <UNK> probabilities
 - Create a fixed lexicon L of size V
 - At text normalization phase, any training word not in L changed to <UNK>
 - Now we train its probabilities like a normal word
 - At decoding time
 - If text input: Use UNK probabilities for any word not in training

Advanced Smoothing Algorithms

- Naïve smoothing algorithms have limited usage and are not very effective. Not frequently used for N-grams.
- However, they can be used in domains where the number of zeros isn't so huge.







Advanced Smoothing Algorithms

- Naïve smoothing algorithms have limited usage and are not very effective. Not frequently used for N-grams.
- However, they can be used in domains where the number of zeros isn't so huge.
- Popular Algorithms:
 - Good-Turing
 - Kneser-Ney

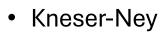




Advanced Smoothing Algorithms

- Naïve smoothing algorithms have limited usage and are not very effective. Not frequently used for N-grams.
- However, they can be used in domains where the number of zeros isn't so huge.
- Popular Algorithms:
 - Good-Turing

Use the count of things we've **seen once** to help estimate the count of things we've **never seen**



LMs: Introduction and Recent Advances





• N_c = Frequency of frequency of c







- N_c = Frequency of frequency of c
- Rohan I am I am Rohan I like to play





- N_c = Frequency of frequency of c
- Rohan I am I am Rohan I like to play

I 3

Rohan 2

Am 2

like 1

to 1

play 1





3

1

1

- N_c = Frequency of frequency of c
- Rohan I am I am Rohan I like to play

Rohan 2

Am 2

like

to 1

play

$$N_1 = 3, N_2 = 2, N_3 = 1$$





- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?

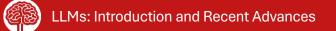




- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18







- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18
- How likely is it that the next bird you see is a new species -- Purple Heron or Painted Stork?





- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18
- How likely is it that the next bird you see is a new species -- Purple Heron or Painted Stork?
 - We will use our estimate of things we saw once to estimate the new things.





- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18
- How likely is it that the next bird you see is a new species -- Purple Heron or Painted Stork?
 - We will use our estimate of things we saw once to estimate the new things.
 - 3/18 (because N₁ = 3)





- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18
- How likely is it that the next bird you see is a new species -- Purple Heron or Painted Stork?
 - We will use our estimate of things we saw once to estimate the new things.
 - 3/18 (because N₁ = 3)
- Assuming so, how likely it is that the new species is Woodpecker?





- You are birdwatching in the Jim Corbett National Park and you have observed the following birds: 10 Flamingos, 3 Kingfishers, 2 Indian Rollers, 1 Woodpecker, 1 Peacock, 1 Crane = 18 birds
- How likely is it that the next bird you see is a woodpecker?
 - 1/18
- How likely is it that the next bird you see is a new species -- Purple Heron or Painted Stork?
 - We will use our estimate of things we saw once to estimate the new things.
 - 3/18 (because N₁ = 3)
- Assuming so, how likely it is that the new species is Woodpecker?
 - Must be less than 1/18





• P^*_{GT} (things with zero frequency) = $\frac{N_1}{N}$







- P^*_{GT} (things with zero frequency) = $\frac{N_1}{N}$
- Unseen (Purple Heron or Painted Stork)
 - C = 0
 - MLE p = 0/18 = 0
 - P^{*}_{GT} (unseen) = N₁/N = 3/18





• P^*_{GT} (things with zero frequency) = $\frac{N_1}{N}$

$$c^* = \frac{(c+1)N_{c+1}}{N_c}$$

- Unseen (Purple Heron or Painted Stork)
 - C = 0
 - MLE p = 0/18 = 0
 - P^{*}_{GT} (unseen) = N₁/N = 3/18





- P^*_{GT} (things with zero frequency) = $\frac{N_1}{N}$
- Unseen (Purple Heron or Painted Stork)
 - C = 0
 - MLE p = 0/18 = 0
 - P^{*}_{GT} (unseen) = N₁/N = 3/18

$$c^* = \frac{(c+1)N_{c+1}}{N_c}$$

- Seen once
 - C = 1
 - MLE p = 1/18

• c* (Woodpecker) = 2 * N₂/N₁
= 2 * 1/3 = 2/3
• P*_{GT} (Woodpecker) =
$$\frac{\frac{2}{3}}{\frac{1}{18}}$$
 = 1/27





Good Turing Estimation

- Numbers from Church and Gale (1991)
- 22 million words of AP Newswire

Count c	Good Turing c*
0	.0000270
1	0.446
2	1.26
3	2.24
4	3.24
5	4.22
6	5.19
7	6.21
8	7.24
9	8.25

Example from Speech and Language Processing book by Daniel Jurafsky and James H. Martin





Good Turing Estimation

- Numbers from Church and Gale (1991)
- 22 million words of AP Newswire

LMs: Introduction and Recent Advances

It looks like $c^* = (c - 0.75)$

Count c	Good Turing c*
0	.0000270
1	0.446
2	1.26
3	2.24
4	3.24
5	4.22
6	5.19
7	6.21
8	7.24
9	8.25

Example from Speech and Language Processing book by Daniel Jurafsky and James H. Martin





Absolute Discounting Interpolation

• Adjusts the probability estimates for n-grams by discounting each count by a fixed amount (usually a small constant) before computing probabilities

$$P_{\text{AbsoluteDiscounting}}(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i) - d}{C(w_{i-1})} + \lambda(w_{i-1})P(w_i)$$







Absolute Discounting Interpolation

• Adjusts the probability estimates for n-grams by discounting each count by a fixed amount (usually a small constant) before computing probabilities

$$P_{\text{AbsoluteDiscounting}}(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i) - d}{C(w_{i-1})} + \lambda(w_{i-1})P(w_i)$$
unigram

• But considering the regular unigram probability has some limitations, as we will see in the upcoming slides.





- My breakfast is incomplete without a cup of ... : coffee/ Angeles?
- Say, in the corpus "Angeles" more prevalent than "coffee"
- However, it is important to note that "Angeles" mostly comes after "Los"
- Instead of regular unigram probability, use **continuation probability**.







- My breakfast is incomplete without a cup of ... : coffee/ Angeles?
- Say, in the corpus "Angeles" more prevalent than "coffee"
- However, it is important to note that "Angeles" mostly comes after "Los"
- Instead of regular unigram probability, use continuation probability.
 - Regular Unigram probability: P(w) : "How likely is w?"
 - P_{continuation}(w): "How likely is w to appear as a novel continuation?"





- My breakfast is incomplete without a cup of ... : coffee/ Angeles?
- Say, in the corpus "Angeles" more prevalent than "coffee"
- However, it is important to note that "Angeles" mostly comes after "Los"
- Instead of regular unigram probability, use **continuation probability**.
 - Regular Unigram probability: P(w) : "How likely is w?"
 - P_{continuation}(w): "How likely is w to appear as a novel continuation?"
- How to compute **continuation probability**?





- My breakfast is incomplete without a cup of ... : coffee/ Angeles?
- Say, in the corpus "Angeles" more prevalent than "coffee"
- However, it is important to note that "Angeles" mostly comes after "Los"
- Instead of regular unigram probability, use continuation probability.
 - Regular Unigram probability: P(w) : "How likely is w?"
 - P_{continuation}(w): "How likely is w to appear as a novel continuation?"
- How to compute continuation probability?
 - Count how many different bigram types each word completes => Normalize by the total number of word bigram types

$$\mathsf{P}_{\text{continuation}}(\mathsf{w}) = \frac{|\{\mathsf{w}_{j-1}, \mathsf{c}(\mathsf{w}_{j-1}, \mathsf{w}) > 0\}|}{|\{(\mathsf{w}_{j-1}, \mathsf{w}_{j}) : \mathsf{c}(\mathsf{w}_{j-1}, \mathsf{w}_{j}) > 0\}|}$$



Intuition: Shannon game

- My breakfast is incomplete without a cup of ... : coffee/ Angeles?
- Say, in the corpus "Angeles" more prevalent than "coffee"
- However, it is important to note that "Angeles" mostly comes after "Los"
- Instead of regular unigram probability, use continuation probability.
 - Regular Unigram probability: P(w) : "How likely is w?"
 - P_{continuation}(w): "How likely is w to appear as a novel continuation?"
- How to compute **continuation probability**?
 - Count how many different bigram types each word completes => Normalize by the total number of word bigram types

 $P_{\text{continuation}}(w) = \frac{|\{w_{j-1}, c(w_{j-1}, w) > 0\}|}{|\{(w_{j-1}, w_j) : c(w_{j-1}, w_j) > 0\}|}$





A common word (Angeles) appearing in only one context (Los) is likely to have a low continuation probability.

Kneser-Ney Smoothing

$$P_{KN}(w_{i} | w_{i-1}) = \frac{\max(c(w_{i-1}, w_{i}) - d, 0)}{c(w_{i-1})} + \lambda(w_{i-1})P_{continuation}(w_{i})$$

where, λ is a normalizing constant (How to define this?)







Kneser-Ney Smoothing

$$P_{KN}(w_{i} | w_{i-1}) = \frac{\max(c(w_{i-1}, w_{i}) - d, 0)}{c(w_{i-1})} + \lambda(w_{i-1})P_{continuation}(w_{i})$$

where, $\boldsymbol{\lambda}$ is a normalizing constant

$$\lambda(w_{i-1}) = \frac{d}{c(w_{i-1})} | \{w : c(w_{i-1}, w) > 0\} |$$







Evaluation of Language Models







Evaluation of a Language Model

• Does our language model prefer good sentences over bad ones?







Evaluation of a Language Model

- Does our language model prefer good sentences over bad ones?
 - Assign higher probability to "real" or "frequently observed" sentences than "ungrammatical" or "rarely observed" sentences
- Terminologies:
 - We optimize the parameters of our model based on data from a training set.
 - We assess the model's performance on unseen **test data** that is disjoint from the training data.
 - An evaluation metric provides a measure of the performance of our model on the test set.







Extrinsic Evaluation

 Measure the effectiveness of a language model by testing their performance on different downstream NLP tasks, such as machine translation, text classification, speech recognition.

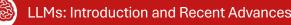






Extrinsic Evaluation

- Measure the effectiveness of a language model by testing their performance on different downstream NLP tasks, such as machine translation, text classification, speech recognition.
- Let us consider two different language models: A and B
 - Select a suitable evaluation metric to assess the performance of the language models based on the chosen task.
 - Obtain the evaluation scores for A and B
 - Compare the evaluation scores for A and B







Intrinsic Evaluation: Perplexity

Intuition: The Shannon Game

- How well can we predict the next word?
 - I always order pizza with cheese and ...
 - The president of India is ...
 - I wrote a ...





Intrinsic Evaluation: Perplexity

Intuition: The Shannon Game

- How well can we predict the next word?
 - I always order pizza with cheese and ...
 - The president of India is ...
 - I wrote a ...
- Observation: The more context we consider, the better the prediction.







Intrinsic Evaluation: Perplexity

Intuition: The Shannon Game

- How well can we predict the next word?
 - I always order pizza with cheese and ...
 - The president of India is ...
 - I wrote a ...
- **Observation:** The more context we consider, the better the prediction.

A better text model is characterized by its ability to assign a higher probability to the correct word in a given context.







The best language model is one that best predicts an unseen test set.

Perplexity is the inverse probability of the test data, normalized by the number of words.

• Given a sentence W consisting of *n* words, the perplexity is calculated as follows:

 $PP(W) = P(w_1 w_2 ... w_n)^{-\frac{1}{n}}$







Thus, for the sentence W, perplexity is:

$$PP(W) = P(w_1 w_2 ... w_n)^{-\frac{1}{n}}$$







Thus, for the sentence W, perplexity is:

$$PP(W) = P(w_1w_2...wn)^{-\frac{1}{n}}$$

Applying Chain Rule:

$$PP(W) = \left(\prod \frac{1}{P(W_{i} | W_{1} W_{2} ... W_{i-1})}\right)^{\frac{1}{n}}$$







Thus, for the sentence W, perplexity is:

$$PP(W) = P(w_1 w_2 ... wn)^{-\frac{1}{n}}$$

Applying Chain Rule:

$$PP(W) = \left(\prod \frac{1}{P(W_{i} | W_{1}W_{2} ... W_{i-1})} \right)^{\frac{1}{n}}$$

Applying Markov Assumption (n = 2), i.e. for bigram LM:

$$\mathsf{PP}(\mathsf{W}) = \left(\prod \frac{1}{\mathsf{P}(\mathsf{w}_{i} | \mathsf{w}_{i-1})}\right)^{\frac{1}{n}}$$







Thus, for the sentence W, perplexity is:

$$PP(W) = P(w_1 w_2 ... wn)^{-\frac{1}{n}}$$

Applying Chain Rule:

$$PP(W) = \left(\prod \frac{1}{P(W_{i} | W_{1}W_{2} ... W_{i-1})} \right)^{\frac{1}{n}}$$

Applying Markov Assumption (n = 2), i.e. for bigram LM:

Minimizing perplexity is the same as maximizing probability.

$$\mathsf{PP}(\mathsf{W}) = \left(\prod \frac{1}{\mathsf{P}(\mathsf{w}_{\mathsf{i}} | \mathsf{w}_{\mathsf{i-1}})}\right)^{\frac{1}{n}}$$







Perplexity and Entropy







- Predicting the next word using a fixed window of previous words.
- Fixed Context Size: Limited to a fixed window of previous words.







- Predicting the next word using a fixed window of previous words.
- Fixed Context Size: Limited to a fixed window of previous words.
- Smoothing techniques address data sparsity.







- Predicting the next word using a fixed window of previous words.
- Fixed Context Size: Limited to a fixed window of previous words.
- Smoothing techniques address data sparsity.
 - But even with smoothing, rare n-grams are hard to predict.







- Predicting the next word using a fixed window of previous words.
- Fixed Context Size: Limited to a fixed window of previous words.
- Smoothing techniques address data sparsity.
 - But even with smoothing, rare n-grams are hard to predict.
- Large vocabulary leads to high memory requirements.





- N-gram LMs suffer from data sparsity and limited context.
 - Predicting the next word using a fixed window of previous words.
 - Fixed Context Size: Limited to a fixed window of previous words.
- Smoothing techniques address data sparsity.
 - But even with smoothing, rare n-grams are hard to predict.
- Large vocabulary leads to high memory requirements.
- High computational cost for large n-grams.



- N-gram LMs suffer from data sparsity and limited context.
 - Predicting the next word using a fixed window of previous words.
 - Fixed Context Size: Limited to a fixed window of previous words.
- Smoothing techniques address data sparsity.
 - But even with smoothing, rare n-grams are hard to predict.
- Large vocabulary leads to high memory requirements.
- High computational cost for large n-grams.
- Lack of generalization to unseen word combinations.







The Need for Richer Representations

Requirements:

- Contextual Understanding: Need for models that understand context beyond fixed windows.
- Semantic Similarity: Ability to capture relationships between words (e.g., synonyms).
- Scalability: Models that can scale to large datasets and handle vast vocabularies efficiently.







Moving to Word Embeddings & Neural LM

In the successive lectures, we will see how representing words (actually, tokens) as vectors and transition to neural LMs solve many of those problems.







Moving to Word Embeddings & Neural LM

In the successive lectures, we will see how representing words (actually, tokens) as vectors and transition to neural LMs solve many of those problems.

- Move from discrete to continuous representations.
- Capture richer semantic information.
- Enable generalization to unseen data.
- Scale to large datasets.





Timeline in Language Modelling

